

Continual Learning

Insights from SplitMNIST

CS 5100 project presentation

Khoury College of Computer Science, Northeastern University

Saurabh Shetty

shetty.sau@northeastern.edu

Siddharth Mittal

mittal.sid@northeastern.edu

TABLE OF CONTENTS



01

The Continual
Learning Problem



02

Approaches
To Continual
Learning



03

Experiments &
Results



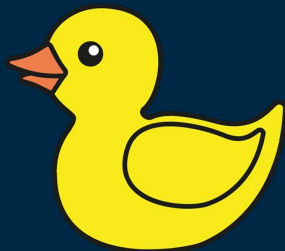
04

Conclusion

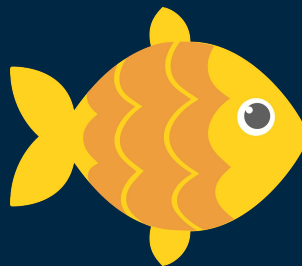


The Continual Learning Problem

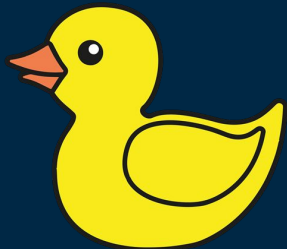
Train a classification model



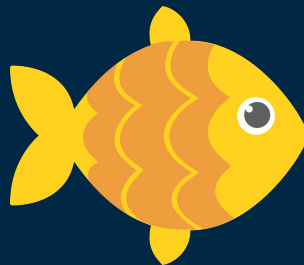
VS



Train a classification model

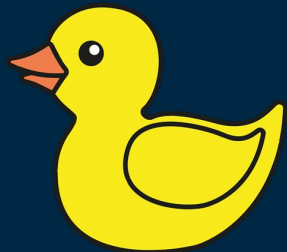


VS

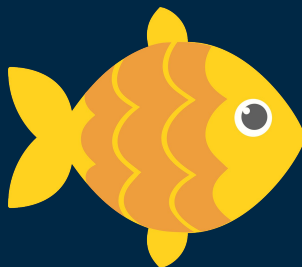


Model 1

Train a classification model



VS

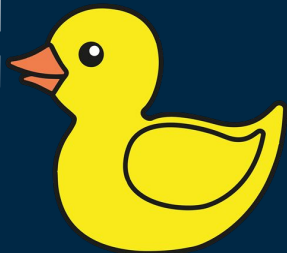


Model 1

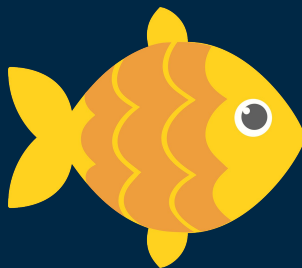
VS



Train a classification model



VS

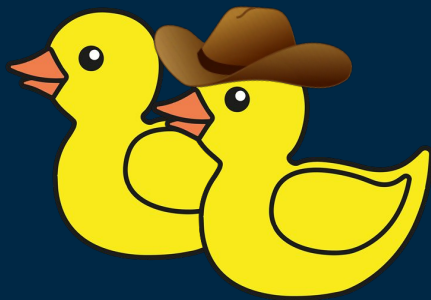


VS

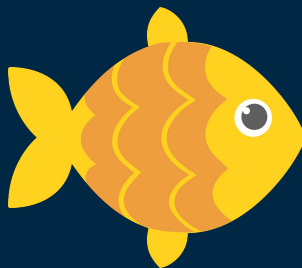


Model 2

Train a classification model



VS



VS



Model 3



How can we improve AI efficiency, scalability and adaptability?

*(Hence making it **sustainable in the long term**)*

Continual Learning



A collection of small squares in various shades of blue, teal, and light blue, scattered in the top right corner of the slide.

Continual Learning

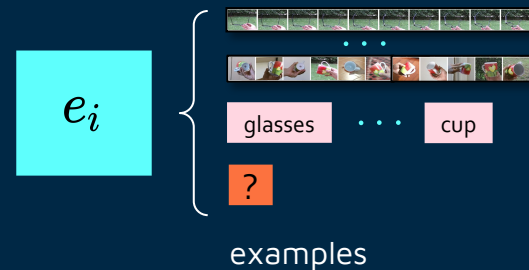
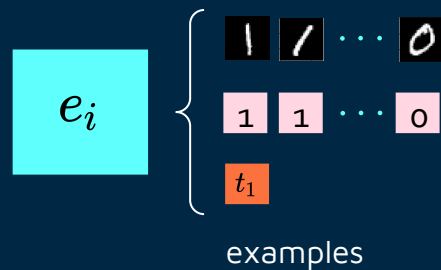
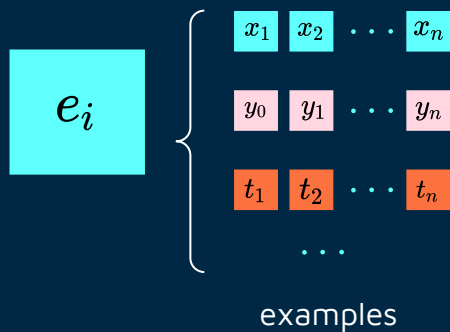
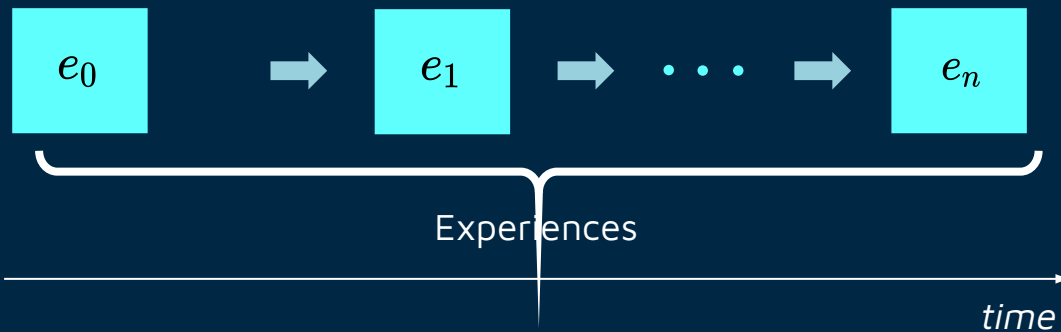
Continual learning is a set of approaches to train machine learning models incrementally, using data samples only once as they arrive.

A small cluster of squares in light blue and teal colors located in the bottom left corner of the slide.

Continual Learning Desiderata

- Higher and **realistic time-scale** where data (and tasks) become available only during time.
- **No access** to previously encountered data.
- **Constant** computational and memory resources (efficiency)
- **Incremental development** of ever more complex knowledge and skills (scalability)
- **Efficiency + Scalability = Sustainability**

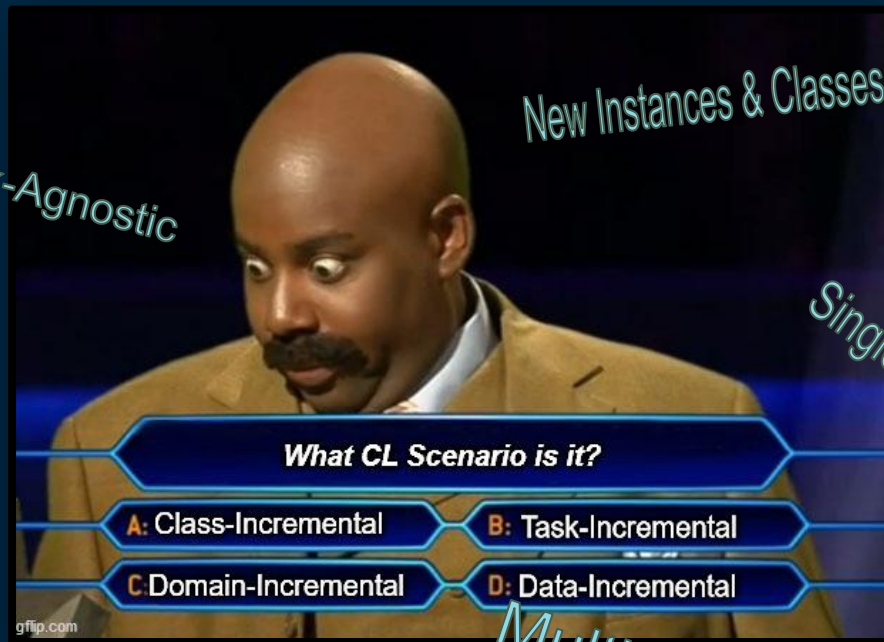
Continual Learning



The background is a dark blue gradient. It features several thin, vertical white lines of varying lengths scattered across the frame. Interspersed among these lines are small squares in three colors: light blue, light orange, and light pink. Some squares are solid, while others are outlined. The overall aesthetic is modern and minimalist.

Continual Learning Scenarios

Task-Agnostic
Task-Free



Multi-Task

Single-Incremental-Task

A Possible Categorization

Experience content type

Task Labels

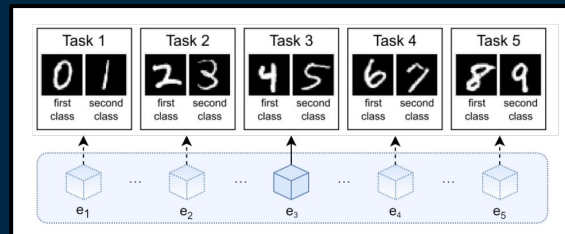
	New Instances (NI)	New Classes (NC)	New Instances and Classes (NIC)
<i>Multi-Task</i>	-	Task Incremental	-
<i>Single-Incremental-Task</i>	Domain-Incremental	Class-Incremental	Data-Incremental
<i>Multiple-Incremental-Task</i>	?	?	?

Our Dataset & Scenario

Mnist Dataset



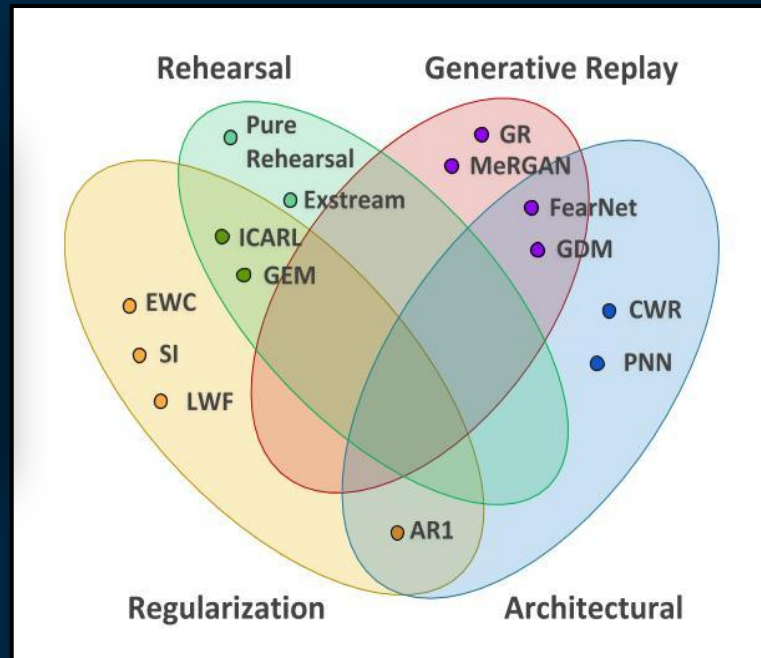
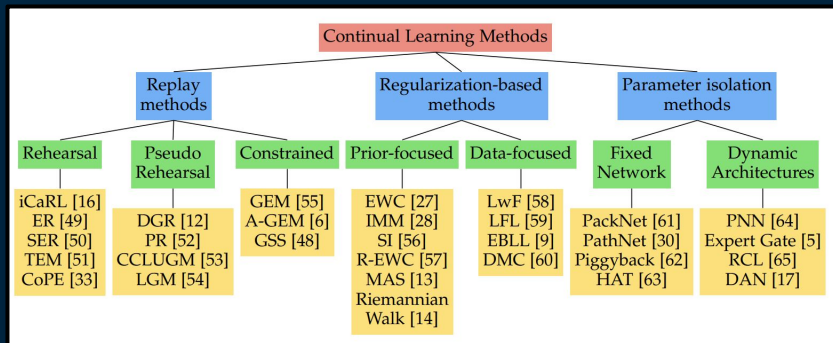
Split Mnist Benchmark





Continual Learning Strategy

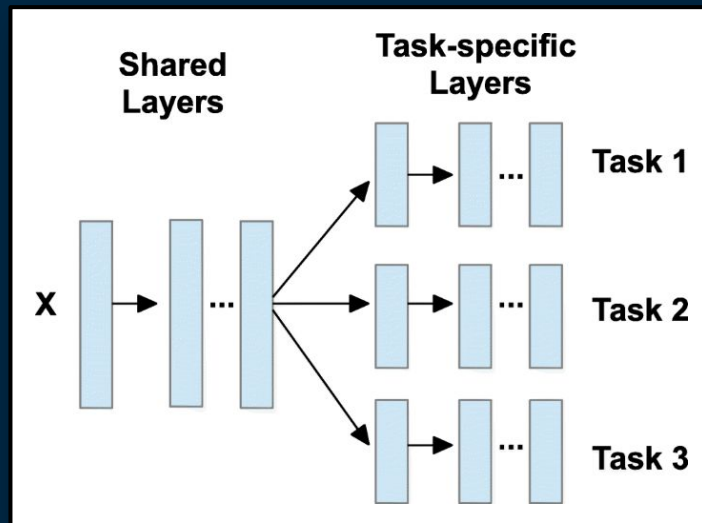
Possible 4-way Categorization



Continual Learning Baselines

Common Baselines / Control Algorithms

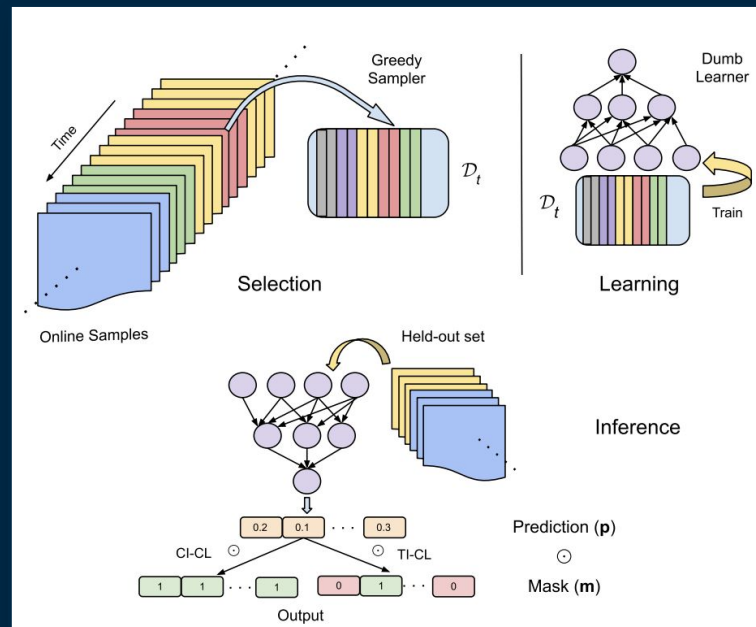
- **Naive / Finetuning** (just continuing backprop)
- **Ensemble**: one model for each experience
- **Cumulative**: for every experience, accumulate all data and re-train from scratch.



GDUMB: Replay based approach

Greedy Sampler and Dumb Learner

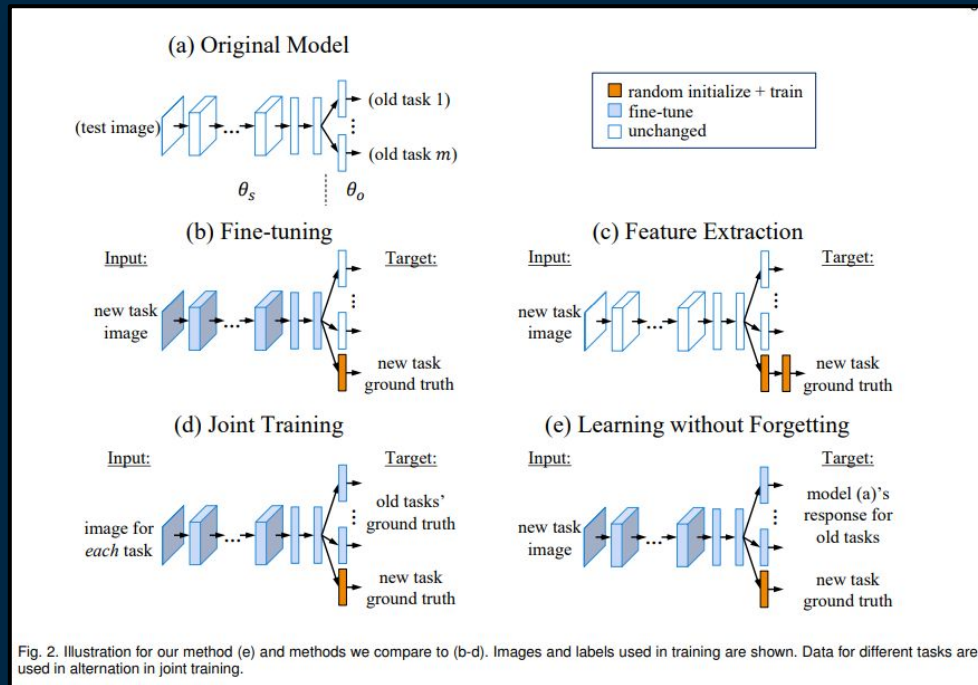
- Interesting paper that sparked **strong discussions** in the CL community
- Note that there's **no knowledge transfer** in this strategy (quite dumb indeed!)
- Despite its simplicity, It was shown to work better than some existing and more complex strategies, **questioning the utility of some benchmarks/metrics** in our field



Learning Without Forgetting (LWF)

Key Aspects

- Straightforward application of **Knowledge Distillation**
- Originally designed for **Task-Incremental settings** can be easily extended to others
- **Efficient single-head implementations** exist
- **Easy to implement** and commonly used



Learning without Forgetting, Li et al, TPAMI 2017

Distilling the knowledge in a neural network, Hinton et al, 2015.

Continuous Learning in Single-Incremental Tasks, Maltoni & Lomonaco, Neural Networks, 2019.

Copy Weights with Re-Init (CWR)

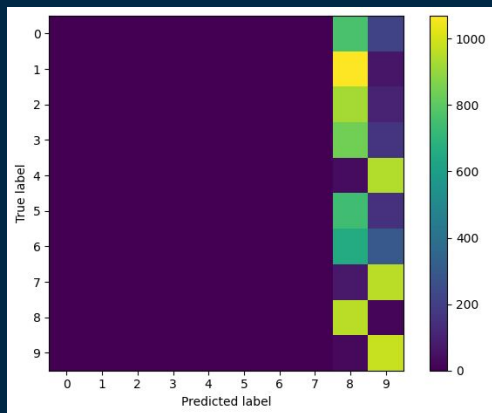
Key Aspects

- Developed for the **fully connected linear classifier** (may be extended to multiple layers)
- **Dual memory system approach**: one for better plasticity, one for memory consolidation
- Very simple and efficient, yet effective solution **agnostic to the experience content (NI, NC, NIC) and specific scenario**

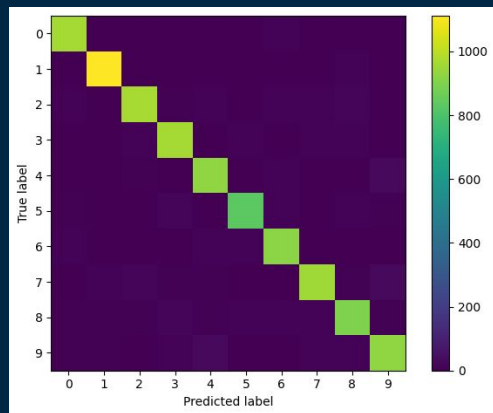
Algorithm 1 CWR* pseudocode: $\bar{\Theta}$ are the class-shared parameters of the representation layers; the notation $cw[j]$ / $tw[j]$ is used to denote the groups of consolidated / temporary weights corresponding to class j . Note that this version continues to work under NC, which is seen here a special case of NIC; in fact, since in NC the classes in the current batch were never encountered before, the step at line 7 loads 0 values for classes in B_i because cw_j were initialized to 0 and in the consolidation step (line 13) $wpast_j$ values are always 0.

```
1: procedure CWR*
2:    $cw = 0$ 
3:    $past = 0$ 
4:   init  $\bar{\Theta}$  random or from pre-trained model (e.g. on ImageNet)
5:   for each training batch  $B_i$ :
6:     expand output layer with neurons for the new classes in  $B_i$ 
       never seen before
7:      $tw[j] = \begin{cases} cw[j], & \text{if class } j \text{ in } B_i \\ 0, & \text{otherwise} \end{cases}$ 
8:     train the model with SGD on the  $s_i$  classes of  $B_i$ :
9:       if  $B_i = B_1$  learn both  $\bar{\Theta}$  and  $tw$ 
10:      else learn  $tw$  while keeping  $\bar{\Theta}$  fixed
11:      for each class  $j$  in  $B_i$ :
12:         $wpast_j = \sqrt{\frac{past_j}{cur_j}}$ , where  $cur_j$  is the number of patterns
          of class  $j$  in  $B_i$ 
13:         $cw[j] = \frac{cw[j] \cdot wpast_j + (tw[j] - avg(tw))}{wpast_j + 1}$ 
14:         $past_j = past_j + cur_j$ 
15:      test the model by using  $\bar{\Theta}$  and  $cw$ 
```

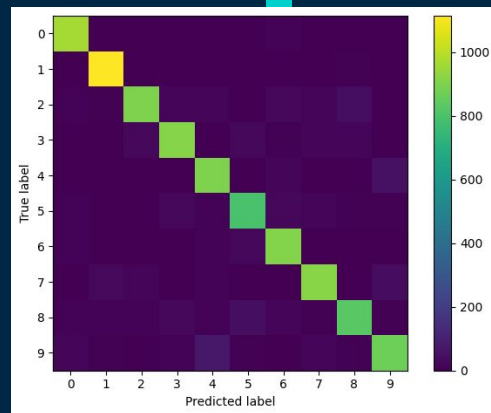
Results (Accuracy Matrix)



Naive

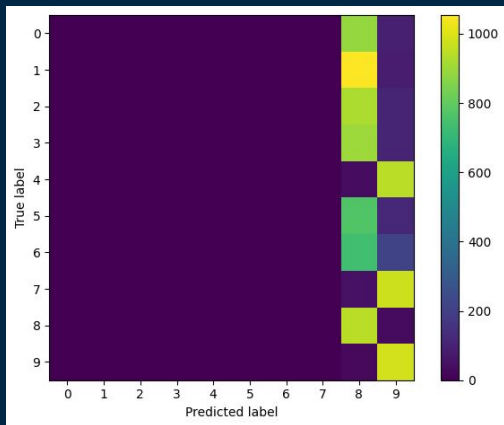


Cumulative

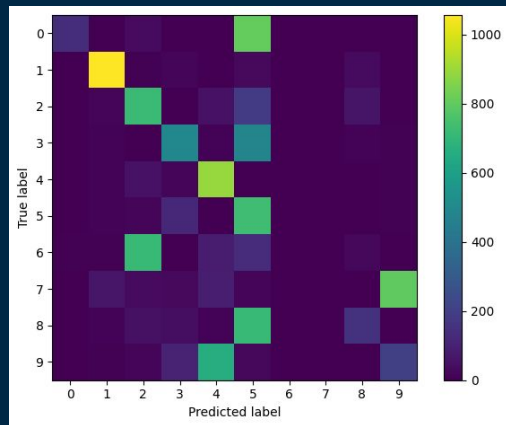


GDumb

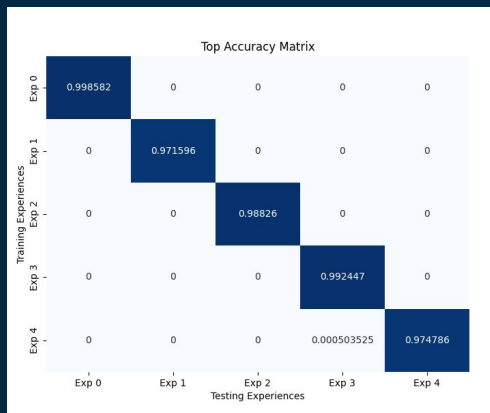
LWF



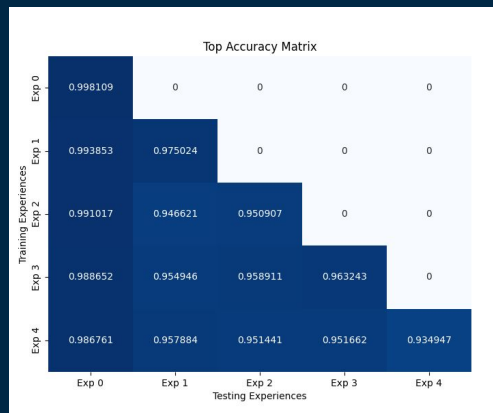
CWR



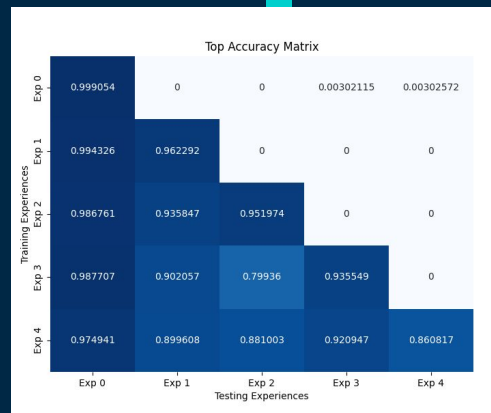
Results (Accuracy)



Naive

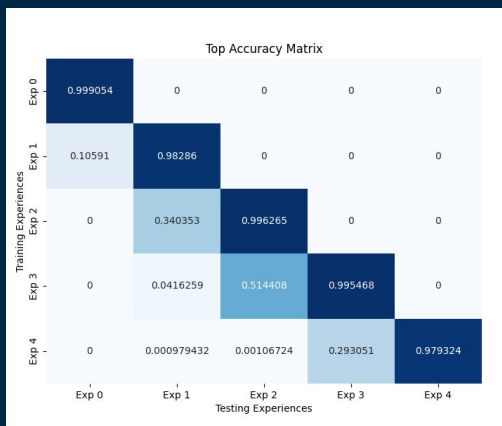


Cumulative

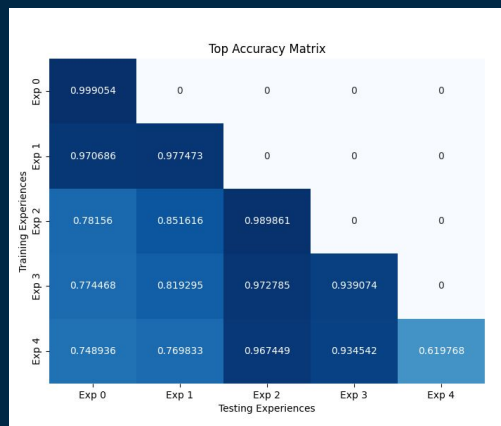


GDumb

LWF



CWR



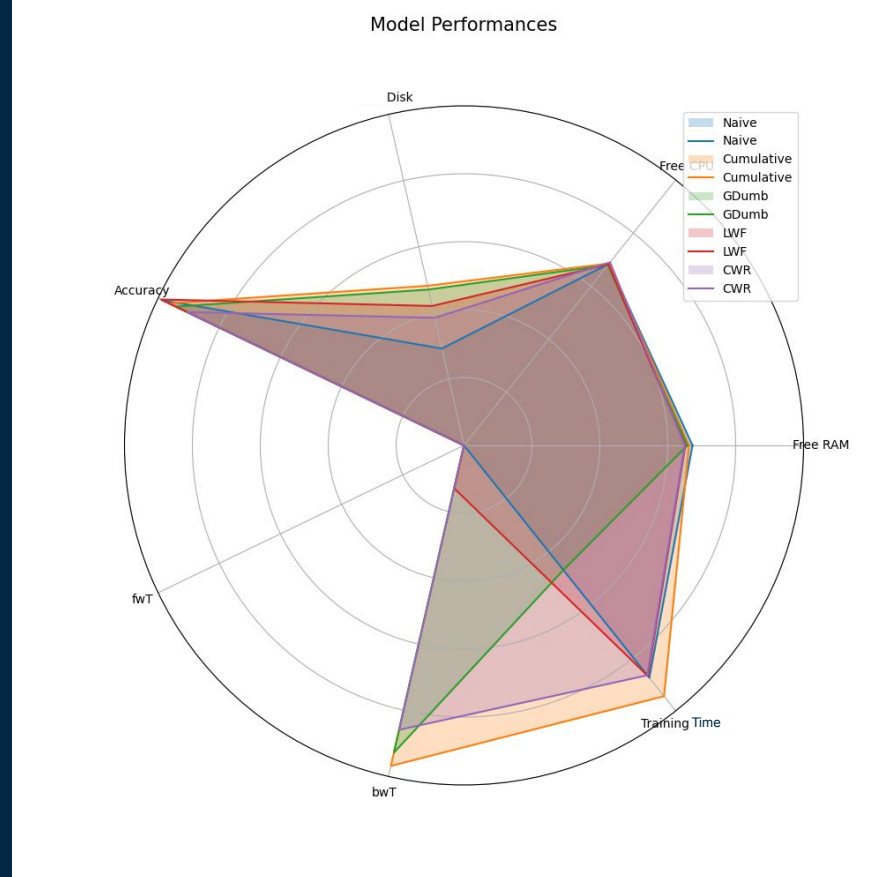
Accuracy metrics

$$\text{Backward Transfer} = \frac{\sum_{i>j}^N \mathcal{R}_{i,j}}{\frac{N(N-1)}{2}}$$

$$\text{Forward Transfer} = \frac{\sum_{i<j}^N \mathcal{R}_{i,j}}{\frac{N(N-1)}{2}}$$

$$\text{In-Domain Accuracy} = \frac{\sum_{i=1}^N \mathcal{R}_{i,i}}{N}$$

Results (Radar Graph)



Future Works

Key Aspects

- Experiment with hybrid models like AR1
- Benchmark on complex datasets like CLEAR (Continual LEArning on Real-World Imagery) dataset
- Explore implementation in robot vision related tasks. I.e. ego-centric video datasets

Q & A